



LAWRENCE  
LIVERMORE  
NATIONAL  
LABORATORY

# Modeling the Office of Science Ten Year Facilities Plan: The PERI Architecture Tiger Team

Bronis R. de Supinski, Sadaf R. Alam, David H. Bailey, Laura Carrington, Chris Daley, Anshu Dubey, Todd Gamblin, Dan Gunter, Paul Hovland, Heike Jagode, Karen Karavanic, Gabriel Marin, John Mellor-Crummey, Shirley Moore, Boyana Norris, Leonid Oliker, Catherine Olschanowsky, Philip C. Roth, Martin Schulz, Sameer Shende, Allan Snaveley, Wyatt Spear, Mustafa Tikir, Jeff Vetter, Pat Worley, Nicholas Wright

May 29, 2009

SciDAC 2009  
San Diego, CA, United States  
June 14, 2009 through June 18, 2009

## **Disclaimer**

---

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

# Modeling the Office of Science Ten Year Facilities Plan: The PERI Architecture Tiger Team

Bronis R. de Supinski<sup>1</sup>, Sadaf Alam<sup>2</sup>, David H. Bailey<sup>3</sup>,  
Laura Carrington<sup>4</sup>, Chris Daley<sup>5</sup>, Anshu Dubey<sup>5</sup>, Todd Gamblin<sup>1</sup>,  
Dan Gunter<sup>3</sup>, Paul Hovland<sup>6</sup>, Heike Jagode<sup>2</sup>, Karen Karavanic<sup>7</sup>,  
Gabriel Marin<sup>8</sup>, John Mellor-Crummey<sup>8</sup>, Shirley Moore<sup>9</sup>,  
Boyana Norris<sup>6</sup>, Leonid Oliker<sup>3</sup>, Catherine Olschanowsky<sup>4</sup>,  
Philip C. Roth<sup>2</sup>, Martin Schulz<sup>1</sup>, Sameer Shende<sup>10</sup>, Allan Snaveley<sup>4</sup>,  
Wyatt Spear<sup>10</sup>, Mustafa Tikir<sup>4</sup>, Jeff Vetter<sup>2</sup>, Pat Worley<sup>2</sup>, and  
Nicholas Wright<sup>4</sup>

<sup>1</sup> Lawrence Livermore National Laboratory, Livermore, California

<sup>2</sup> Oak Ridge National Laboratory, Oak Ridge, Tennessee

<sup>3</sup> Lawrence Berkeley National Laboratory, Berkeley, California

<sup>4</sup> San Diego Supercomputing Center, San Diego, California

<sup>5</sup> University of Chicago, Chicago, Illinois

<sup>6</sup> Argonne National Laboratory, Argonne, Illinois

<sup>7</sup> Portland State University, Portland, Oregon

<sup>8</sup> Rice University, Houston, Texas

<sup>9</sup> University of Tennessee - Knoxville, Knoxville, Tennessee

<sup>10</sup> University of Oregon, Eugene, Oregon

E-mail: bronis@llnl.gov, alamsr@ornl.gov, dhbailey@lbl.gov, lcarrington@sdsc.edu,  
cdaley@flash.uchicago.edu, a-dubey1@uchicago.edu, tgamblin@llnl.gov,  
dkgunter@lbl.gov, hovland@mcs.anl.gov, jagode@eecs.utk.edu, karavan@cs.pdx.edu,  
mgabi@cs.rice.edu, johnmc@cs.rice.edu, shirley@cs.utk.edu, norris@mcs.anl.gov,  
loliker@lbl.gov, cmills@sdsc.edu, rothpc@ornl.gov, schulzm@llnl.gov,  
sameer@cs.uoregon.edu, allans@sdsc.edu, wspear@cs.uoregon.edu, mtikir@sdsc.edu,  
vetter@ornl.gov, worleyph@ornl.gov, nwright@sdsc.edu

**Abstract.** The Performance Engineering Institute (PERI) originally proposed a tiger team activity as a mechanism to target significant effort to the optimization of key Office of Science applications, a model that was successfully realized with the assistance of two JOULE metric teams. However, the Office of Science requested a new focus beginning in 2008: assistance in forming its ten year facilities plan. To meet this request, PERI formed the Architecture Tiger Team, which is modeling the performance of key science applications on future architectures, with S3D, FLASH and GTC chosen as the first application targets. In this activity, we have measured the performance of these applications on current systems in order to understand their baseline performance and to ensure that our modeling activity focuses on the right versions and inputs of the applications. We have applied a variety of modeling techniques to anticipate the performance of these applications on a range of anticipated systems. While our initial findings predict that Office of Science applications will continue to perform well on future machines from major hardware vendors, we have also encountered several areas in which we must extend our modeling techniques in order to fulfill our mission accurately and completely. In addition, we anticipate that models of a wider range of applications will reveal critical differences between expected future systems, thus providing guidance for future Office of Science procurement decisions, and will enable DOE applications to exploit machines in future facilities fully.

## 1. Introduction

Sustained performance improvements are integral to the DOE Office of Science SciDAC program’s mission to advance large-scale scientific modeling and simulation. Simulation is a key investigative technique for disciplines where experimentation is expensive, dangerous, or impossible. Increased performance can enable faster simulations and more timely predictions, or it can be used to increase the accuracy of existing physical models, enabling more predictive simulations. Research enabled by the SciDAC program will have far-reaching effects in fields such as basic energy, biology, environmental science, fusion energy, and high-energy physics.

The Performance Engineering Research Institute (PERI) tiger team activity targets critical SciDAC performance needs. The original intent was for each tiger team to focus the efforts of several PERI researcher on improving performance of an Office of Science application, with the application selected based on Office of Science mission objectives and application readiness for the focused effort. Thus, each tiger team was envisioned as a relatively short-term activity (six months to at most one year). In 2007, our tiger teams had a positive impact on two key DOE applications participating in the JOULE metric. We improved the performance of a turbulent combustion code (S3D [1]) on Oak Ridge’s Cray XT5 Jaguar system by 13%. Similarly, we improved the performance of the Gyrokinetic Toroidal Code (GTC) [2, 3] by 10% on Jaguar and by 15% on Argonne National Laboratory’s (ANL’s) Intrepid Blue Gene/P system.

In 2008, the Office of Science requested that PERI provide assistance in its ten year facilities plan. In particular, they wanted PERI to provide guidance in how key applications would perform across the range of future systems expected to be offered by major vendors in that period. Thus, we redefined the scope of the tiger team activity to handle this request and started the PERI Architecture Tiger Team. This team’s goal is to model the behavior of selected applications and to predict their performance on anticipated future systems instead of to improve their performance on current systems. For this activity to fulfill the request, we must consider a wider range of Office of Science applications, and evaluate the suitability of current and future HPC architectures for the applications. This broader scope has led us to include nearly all PERI researchers on the Architecture Tiger Team.

Several factors complicate the Architecture Tiger Team’s charge. Large scale simulations are complex software artifacts for which the performance depends on input and frequently evolves during the course of a simulation. Further, small source code changes can lead to significant performance changes. Thus, modeling their performance across a variety of existing architectures remains a topic of research. For example, modeling at larger scales than are currently run requires changes to most existing modeling methodologies. In order to model the performance for systems that will emerge over the ten year period, we must not only overcome these challenges but also anticipate how the software, as well as the hardware, will evolve.

For these reasons, we have developed a three-part, iterative plan, with each iteration focusing our modeling effort on a different (or growing) set of applications. First, we extensively measure the performance of the applications at scale with a variety of state-of-the-art performance analysis tools. These measurements ensure that we have appropriate versions of the applications: although we are no longer focused on optimization, we still apply our expertise in this direction. This activity not only ensures that we base our models on an appropriate version but can also provide some benefit to the application teams.

Second, we use these measurements and other data to create predictive performance models that estimate the scaling properties of current applications on future hardware. In the first iteration of the Architecture Tiger Team, we have applied this strategy to three Office of Science early science applications: S3D, GTC, and FLASH, an astrophysical thermonuclear flash simulation. In this paper, we detail the preliminary results of this study, which indicates that these applications will perform well across the breadth of anticipated architectures.

In the third part of our process, we report findings to the Office of Science and work with

them to select the applications for the next iteration. We are currently engaged in that selection process for the Architecture Tiger Team’s second iteration. We are employing criteria that both reflect the importance of the applications to the Office of Science’s mission and attempt to capture the breadth of characteristics of its applications. Simply put, we must ensure that the ten year facilities plan reflects the range of needs of the Office of Science’s broad mission.

The rest of this report is organized as follows. We summarize key tools for our large-scale performance measurement activity in Section 2. Section 3 describes our performance modeling techniques. In Section 4, Section 5 and Section 6, we detail our initial findings with the S3D, FLASH, and GTC codes. We then state our initial conclusions and lessons learned for this on-going activity, including guidance in selecting the next set of applications, in Section 7.

## 2. Measurement

We have used a wide variety of performance analysis tools to characterize the behavior of S3D, GTC and FLASH on current Office of Science platforms at scale. The large volumes of performance data complicate performance measurement on systems such as Jaguar and Intrp Reid. Because these modern parallel applications can have dynamic behavior, understanding their performance can potentially require measuring all application processes. However, for large machines, the overhead of data collection and aggregation can perturb running applications, making the measurements and, thus, the models that we derive from them, inaccurate. Further, too much performance data can make analysis prohibitively expensive.

### 2.1. Performance Analysis Tools

To address these challenges, we have employed a wide variety of tools for measuring performance data of the three applications that we studied in the first iteration of the Architecture Tiger Team. We briefly describe some of our key performance tools in this section; we present results of applying them in Sections 4, 5, and 6.

*2.1.1. Vampir* Vampir [4, 5], an MPI tracing tool developed at TU Dresden, records timelines of all MPI events in a large application run. Vampir stores the collected data to disk for postmortem analysis using a sophisticated set of tools. Users can view summary data for any segment of a trace, and they can also view the trace at varying levels of detail using Vampir’s zoom features. Vampir traces can consume large amounts of space on disk, but Vampir provides parallel analysis tools to ease some of the burden of analyzing this data.

*2.1.2. mpiP* mpiP [6], an MPI profiling tool, measures cumulative time spent in all MPI call sites across all processes in an application. Like other profiling tools, mpiP only collects statistical information, as opposed to full trace data like Vampir. mpiP generates a single file, which is much smaller than a full trace file, but which loses timing information.

*2.1.3. TAU* The TAU [7] suite of parallel performance tools combines the tracing functionality of Vampir and the profiling capabilities of mpiP with sophisticated facilities for source instrumentation as well as facilities for detailed analysis of performance data. While Vampir and mpiP only support measurement of MPI events, TAU can measure specific functions, code regions, and user-defined events in parallel applications. The user must recompile his or her application with the TAU compilers and then re-run a parallel job. TAU then outputs a trace or profile as desired. TAU also provides extensive data mining and analysis tools for processing information after it has been measured and stored.

*2.1.4. Libra* Libra [8] is a tool for scalable load-balance analysis developed at Lawrence Livermore National Laboratory and the University of North Carolina at Chapel Hill. Unlike full trace tools, Libra uses aggressive, lossy wavelet compression to reduce the volume of load-balance data significantly before recording it. Libra can achieve 100:1 to 1000:1 compression on load-balance data, and it also provides a scalable client-side visualization tool for viewing recorded traces. Libra records measured code regions by call site.

## *2.2. Platforms*

In the first iteration of the Architecture Tiger Team, we have conducted extensive measurements of our target applications' performance on two leadership-class systems. The first system is Argonne National Laboratory's Blue Gene/P system, Intrepid. Intrepid contains 163,840 PowerPC 450 cores running at 800 Mhz, and sustained LINPACK performance of 450 Teraflops. The second system is the Cray XT4 Jaguar system at Oak Ridge National Laboratory. Jaguar contains 31,328 Opteron cores running at 2.1 Ghz, and has sustained performance of XXX Teraflops. Both systems use quad-core nodes.

Intrepid and Jaguar have slightly different network and I/O configurations. Jaguar uses a 3D mesh network for communication between nodes, while Intrepid uses a full 3D torus and also uses a tree network and a barrier network for collective communication. Both systems have dedicated I/O nodes that relay I/O operations between applications and the parallel filesystem. On Intrepid, I/O nodes are internal nodes in the tree network. Compute nodes communicate with I/O nodes through the tree network, and the I/O nodes communicate with the parallel filesystem over Myrinet links. On Jaguar, the I/O nodes are situated along one side of the 3D mesh, and compute nodes communicate with them over the mesh network.

## **3. Modeling**

We have developed several techniques to predict performance of DOE applications on leadership class facilities. In this section, we give a brief overview of these techniques, while we discuss the results of applying them to S3D, FLASH and GTC in Sections 4, 5 and 6.

### *3.1. Convolving Machine Profiles with Application Signatures*

To predict the performance of applications on future architectures, we have developed an approach [9] that separates application-specific measurements from machine-specific measurements. Our approach involves two key components:

**Machines Profiles** that characterize the rates at which a machine can (or is expected to) carry out fundamental operations abstract from any particular application;

**Application Signatures** that characterize the fundamental operations that an application must execute independent of any particular machine.

Our approach enables performance predictions of applications on current systems by convolving application signatures with profiles of the existing systems, and on future systems by convolving the application profiles with profiles generated from the *expected* performance parameters of the future systems. Conceptually, a *convolution* defines an algebraic mapping of application signatures onto runtimes to arrive at a performance prediction.

Given an application profile  $A$ , and a machine profile  $M$ , we define  $P$ , a matrix of runtimes, such that  $p_{ij} = \sum_{k=1}^p a_{ik}m_{kj}$ , or:

$$\begin{bmatrix} p_{11} & p_{12} & p_{13} \\ p_{21} & p_{22} & p_{23} \\ p_{31} & p_{32} & p_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} & a_{15} \\ a_{21} & a_{22} & a_{23} & a_{24} & a_{25} \\ a_{31} & a_{32} & a_{33} & a_{34} & a_{35} \end{bmatrix} \begin{bmatrix} m_{11} & m_{12} & m_{13} \\ m_{21} & m_{22} & m_{23} \\ m_{31} & m_{32} & m_{33} \\ m_{41} & m_{42} & m_{43} \\ m_{51} & m_{52} & m_{53} \end{bmatrix}$$

$$e.g., p_{32} = a_{31}m_{12} + a_{33}m_{22} + a_{33}m_{32} + a_{34}m_{42} + a_{35}m_{53}$$

The rows of  $P$  correspond to applications while the columns correspond to systems and each  $p_{ij}$  is the expected runtime of application  $i$  on system  $j$ . The rows of  $A$  are applications, and the columns are columns are operation counts. Any row of  $A$  is the *signature* for application  $i$ . Likewise, rows of  $M$  are bandwidths measured by some benchmark for a particular system while each column is the *profile* of a particular system.

This approach is generic, and we could apply it to measurements from any of the tools mentioned in Section 2 to produce application signatures targeted at particular types of analysis. However, in order to reflect the impact of timing considerations, we currently use traces of memory operations to characterize the fundamental operations of computational code regions and message traces similar to those produced by Vampir [4, 5]. We then use cache simulation to convolve the signatures of the computational regions with characterizations of the memory system obtained with the MultiMAPS benchmark from SDSC. Finally, we use a high-level network simulation such as Dimemas [10, 11] or SDSC’s PSiNS [12] to convolve message trace signature with simple network signatures that capture latency and bandwidth and the predicted performance of the computational regions.

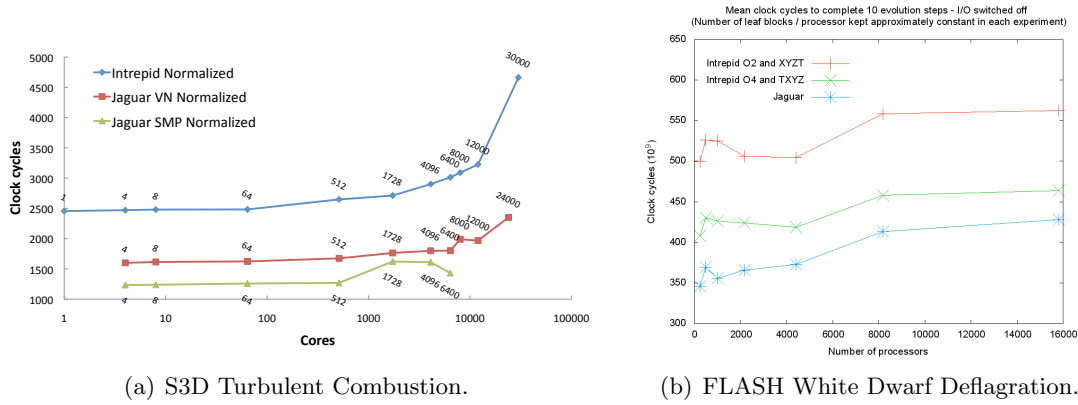
### 3.2. Modeling Assertions

An alternative modeling strategy, called *Modeling Assertions* (MA) [13], constructs symbolic models of application performance. This technique allows users to annotate their code with expressions revealing the relationships among important input parameters, computation, and communication. These annotations, in the form of pragmas or directives, capture the anticipated performance, in terms of time or other metrics, such as cache misses or floating point operations. As the application runs, the MA library checks the model against application structure and key model input parameters. Symbolic performance models complement empirically derived models because the symbolic models expose sensitivities across important parameters and can be scaled to any parameter range.

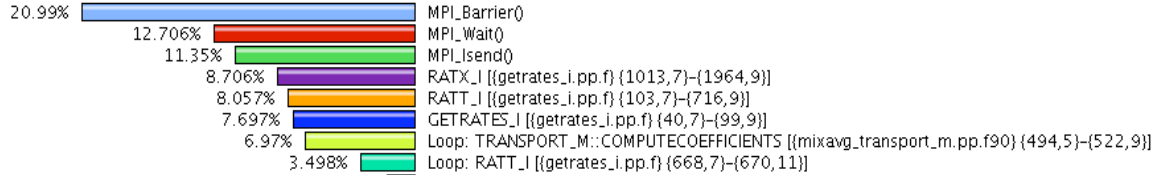
### 3.3. Load Balance Modeling

In addition to models for application runtimes, we have developed models of the load balance properties of large systems. Libra’s compressed representation of system-wide load balance traces uses a wavelet approximation to allow for multiscale representations of load balance properties. The structure of this approximation supports extraction of a low-resolution model of an application’s load balance without recording exhaustive measurements.

This type of compact model will enable us to verify, possibly with distributed extensions of MA, that specific processes have workloads within their expected bounds. Further, our low-resolution models will eventually enable *prediction*, within some confidence, the evolution of load on large systems, and to incorporate dynamically derived load balance guidelines to application-specific load balance components. This solution will dramatically reduce the burden of large-scale measurement on the application developers, enabling them to concentrate on how best to redistribute work within their applications.



**Figure 1.** Normalized Weak Scaling for Office of Science Codes on Intrepid and Jaguar.



**Figure 2.** Most Time Consuming S3D Routines on Jaguar at Scale

#### 4. S3D

S3D [1] is the state of the art turbulent combustion simulation. The code, which was developed at the Combustion Research Facility at Sandia National Laboratory in Livermore, California, won a 2007 INCITE award for six million hours on the XT3/4 Jaguar system at ORNL’s National Center for Computational Sciences. S3D solves the compressible reacting Navier-Stokes equations using high-fidelity numerical methods. Principal components include an eighth-order finite-difference solver, a fourth-order Runge-Kutta integrator, a hierarchy of molecular transport models and detailed chemistry. The use of Direct Numerical Simulation (DNS) enables scientists to study the microphysics of turbulent reacting flows, as this gives full access to time-resolved fields and provides physical insight into chemistry turbulence interactions. Perhaps more importantly, S3D is critical for accurate simulations of larger systems. The detail afforded by the DNS model enables the development of reduced model descriptions that can be used in macro-scale simulations of engineering-level systems.

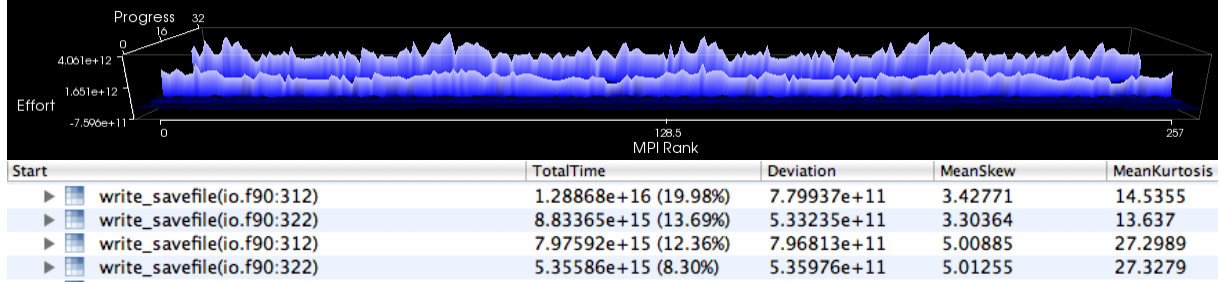
S3D is architected for scalability. It uses a 3D domain decomposition, where each MPI process manages an equal number of grid points and has the same computational load. Interprocessor communication in this decomposition is only between nearest neighbors, and S3D uses large messages and can overlap communication and computation. All-to-all communication is only required for monitoring and synchronization ahead of I/O.

##### 4.1. Measurement of S3D

Figure 1(a) shows our measurements of the weak scaling behavior of S3D on Intrepid and Jaguar. On both systems, S3D scales almost perfectly up to 4,096 cores. After this point, runtimes begin to increase, until at 30,000 processes the runtime is twice that of the baseline, 4-core run on Intrepid. On Jaguar, our 24,000-core run took approximately 30% longer than the baseline run.

We used optimized TAU instrumentation in order to determine the cause of S3D runtime increases at scale. Figure 2 shows the top eight entries in the profile; the first three are





**Figure 3.** Libra MPI Load Balance Profile: Two Checkpoints on Intrepid at 16,384 Cores

(`MPI_Barrier()`, `MPI_Wait()`, and `MPI_Isend()`). S3D spends the bulk of its time in these routines at large core counts with its default IO scheme that writes a file per MPI task. The next most time consuming routines (three subroutines: `RATX_I`, `RATT_I`, and `GETRATES_I`; and two loops) are in the parallel solver. Further correlation studies with TAU showed that the performance of `MPI_Barrier()`, `MPI_Wait()` were the cause of the scaling problems.

Results applying Libra to S3D reveal the underlying issue: the default IO configuration taxes the IO systems excessively. Figure 3 shows a Libra plot of two different checkpoint operations with MPI ranks on the X-axis and time for each task in each checkpoint on the Z-axis. Clearly, the variance in the MPI times shown in Figure 2 reflect a load imbalance caused by highly variable times to complete the IO phase across the MPI tasks. Other IO configurations, including one that performs writes from a subset of the MPI tasks, offer better scaling performance. We are currently working with Petascale Data Storage Institute (PDSI) to understand and to model S3D IO behavior as it is a critical component in its overall performance.

#### 4.2. Modeling of S3D

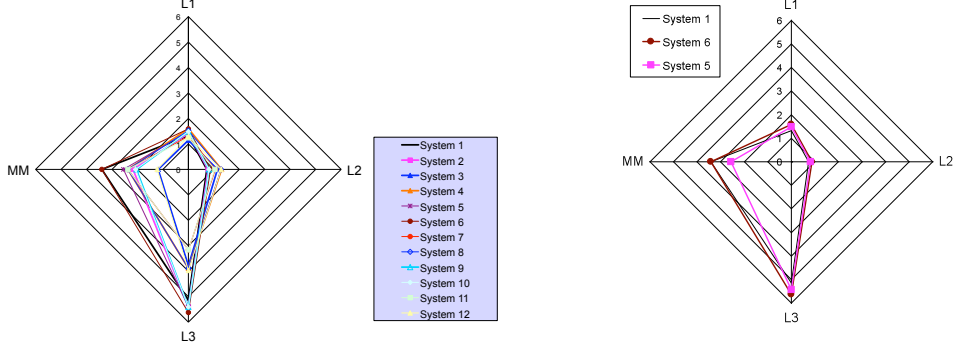
We now detail performance models of S3D’s computational regions, leaving models that include its IO for future work. The Kiviat diagram in Figure 4(a) shows anticipated memory system parameters for several hardware vendors, which we anonymize here due to NDA considerations. The four axes are memory bandwidths for L1, L2, and L3 caches and for main memory (MM). System 1 represents these parameters for Jaguar. In the diagram, the noticeable difference between current and future systems are the significant changes in L3 and main memory bandwidth. Our modeling analysis explores the impact that this difference will have on S3D.

Table 1 shows the results of convolving the machine profiles shown in Figure 4(a) with S3D memory profiles. These results indicate that the differences in memory system performance *will* impact S3D runtimes significantly. We predict that S3D’s  $C_2H_4$  problem will perform well on all expected future systems but it will perform best on those systems with the most main memory bandwidth. Although not shown here, we note that this effect is not true for all applications – for example, our predictions indicate the memory system differences will provide little benefit to WRF, a weather forecasting simulation.

An important consideration for models of S3D is that its memory behavior is scale invariant. We have compared memory traces across a range of job sizes and found that they are very consistent under weak scaling. Thus, we have previously shown on Jaguar that we can predict S3D performance as we scale the number of MPI tasks directly from traces of smaller runs. Thus, we expect the results shown in Table 1 to hold for much larger systems.

## 5. FLASH

FLASH is a parallel, block-structured AMR code designed for compressible reactive flows. Its capabilities span a broad range of applications, from laser-driven shock instabilities to fusion



(a) Twelve Machine Profiles Used with S3D

(b) Three Machine Profiles Used with FLASH

**Figure 4.** Memory profiles used to predict code performance for machines out to 2012.

CPUs	System 1	System 2	System 3	System 4	System 5	System 6	System 7	System 8
8	1	.72	.78	.66	.77	.81	.74	.61
64	1	.73	.78	.67	.77	.81	.73	.62
512	1	.72	.78	.66	.77	.81	.74	.62

**Table 1.** Prediction of S3D  $C_2H_4$  Benchmark Performance on Systems Anticipated by 2012

burn in type Ia supernovae. FLASH has demonstrated wide portability by running successfully on many leadership class systems. It is fully modular, in that its components can be used to create many different astrophysical applications.

### 5.1. Measuring FLASH

We measure FLASH's performance on leadership systems as the number of CPU cycles to complete 10 FLASH time steps with I/O disabled. Figure 1(b) shows that FLASH scales well for a white dwarf deflagration simulation on both Intrepid and on Jaguar. We observe that the curves are similar, increasing slightly as the MPI task count increases. We note that using a slightly different assignment of MPI tasks to processors makes a significant performance difference on Intrepid, with the normalized performance very similar to that on Jaguar, with indications that they might cross at even higher core counts.

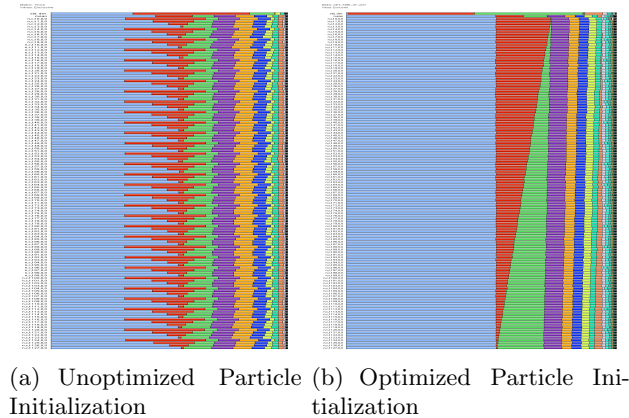
### 5.2. Modeling FLASH

As with S3D, we modeled the memory behavior of the computational phases of FLASH for three anonymous future architectures for which Figure 4(b) shows memory system profiles. Our predictions in Table 2 show preliminary results for 128, 256 and 384 cores on those systems as well as Jaguar and Lonestar, the Sun Infiniband cluster at the Texas Advanced Computing Center. The reasonable accuracy on the existing systems lends confidence to the predictions on future systems. Overall, the results demonstrate that FLASH also will perform well on anticipated future systems and will benefit from improvements to main memory bandwidth.

Our study of FLASH memory traces found that they are not scale invariant. Thus, although we expect FLASH to scale well based on our empirical measurements, we need to extend our modeling techniques to extrapolate memory traces from a set of traces gathered from smaller runs. We are currently pursuing this research direction and initial predictions based on this

CPUs	Lonestar			Jaguar			Prediction of Systems		
	Predicted	Real	% err	Predicted	Real	% err	Sys. 1	Sys. 5	Sys. 6
128	246	227	8.2	285	258	10.3	195	188	148
256	127	131	-2.8	145	138	5.5	99	95	75
384	86	103	-16.4	99	97	2.0	67	64	51

**Table 2.** Performance Predictions for FLASH on Current and Future Systems



**Figure 5.** GTC Load Balance Without and With Optimized Initialization

work confirm that FLASH will scale reasonably well on future systems although additional work remains to confirm this hypothesis on the very large processor counts anticipated during the horizon of the Office of Science’s ten year plan.

## 6. GTC

The Gyrokinetic Toroidal Code (GTC) is a particle-in-cell code to study microturbulence in magnetically confined fusion plasmas. GTC solves the gyro-averaged Vlasov equation and the Gyrokinetic Poisson equation. This global code simulates the entire torus rather than just a flux tube. Written in Fortran 90/95, GTC was originally optimized for superscalar processors, but is now a massively parallel code and frequently uses 1024 or more cores.

Our results with GTC are preliminary, but they demonstrate the value of our measurement activity in ensuring that we use a valid version for modeling. Our initial version of GTC had portions of the particle initialization commented out. Although it still executed correctly, this change led to a significant load imbalance even at small scales. Figure 5(a) shows the TAU profile for GTC on 128 cores of Jaguar with the unoptimized particle initialization. A different color represents each routine in the figure and each MPI task is a row in the profile. The figure clearly shows the load imbalance in the staggering of some routines based on the per-task workload. The optimized particle initialization corrects this imbalance, as the profile in Figure 5(b) shows. The optimizations result in the profile bars for the routines lining up much more evenly across tasks, thus improving GTC’s runtime. This test demonstrates the importance of proper test code configuration for performance modeling.

## 7. Conclusion

We formed the PERI Architecture Tiger Team in order to assist the Office of Science in formulating its ten year facilities plan. Our role is to provide confidence that future leadership class systems will serve the broad range of simulations needed for the Office of Science to fulfill

its mission. The first iteration of our iterative, three phase plan is nearing completion. Its results clearly demonstrate that S3D will perform well on anticipated future platforms with a preference for those that provide the highest main memory bandwidth. Our initial models for FLASH provide similar expectations although we must complete additional research that will enable scaling models for applications that do not exhibit scale-invariant memory reference behavior.

We are currently completing work on measuring and modeling GTC. Our measurement activity of both GTC and S3D demonstrated its value by ensuring that we model an appropriate version of the software. As we complete this first iteration, we are preparing for the next. We are recommending to the Office of Science that we select the next set of applications with a careful eye towards those that exhibit performance differences on current platforms and that stress different aspects of memory and network performance. Thus, we anticipate focusing on at least one latency sensitive application and one bandwidth sensitive application.

## Acknowledgments

This research is sponsored by the Office of Advanced Scientific Computing Research; U.S. Department of Energy. The work of de Supinski, Gamblin and Schulz was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under contract DE-AC52-07NA27344 (LLNL-CONF-xxxxxx). The work of Alam, Jagode, Roth, Vetter and Worley was performed at the Oak Ridge National Laboratory, which is managed by UT-Battelle, LLC under Contract No. DE-AC05-00OR22725. Accordingly, the U.S. Government retains a nonexclusive, royalty free license to publish or reproduce the published form of this contribution, or allow others to do so, for U.S. Government purposes.

This research used resources of the National Center for Computational Sciences at Oak Ridge National Laboratory, which is supported by the Office of Science of the Department of Energy under Contract DE-AC05-00OR22725, and resources of the Argonne Leadership Computing Facility at Argonne National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under contract DE-AC02-06CH11357.

## References

- [1] Hawkes E R and Chen J H 2004 *Combustion and Flame* **138** 242–258
- [2] Lee W W 1983 *Physics of Fluids* **26** 556
- [3] Lee W W 1987 *Journal of Computational Physics* **72** 243
- [4] Brunst H, Hoppe H C, Nagel W E and Winkler M 2001 Performance optimization for large scale computing: The scalable VAMPIR approach *Proceedings of the 2001 International Conference on Computational Science (ICCS 2001)* (San Francisco, CA) pp 751–760
- [5] Brunst H, Kränzlmüller D and Nagel W 2005 *The International Series in Engineering and Computer Science, Distributed and Parallel Systems* **777** 92–102
- [6] Vetter J and Chambreau C 2005 mpiP: Lightweight, scalable mpi profiling URL <http://www.llnl.gov/CASC/mpip>
- [7] Shende S and Maloney A 2006 *International Journal of High Performance Computing Applications* **20** 287–331
- [8] Gamblin T, de Supinski B R, Schulz M, Fowler R J and Reed D A 2008 Scalable load-balance measurement for SPMD codes *Supercomputing 2008 (SC'08)* (Austin, Texas) pp 46–57
- [9] Snavely A, Wolter N and Carrington L 2001 Modeling application performance by convolving machine signatures with application profiles *IEEE Workshop on Workload Characterization, 2001.*
- [10] Labarta J, Girona S and Cortes T 1997 *Parallel Computing* **23** 23–34
- [11] Girona S, Labarta J and Badia R M 2000 Validation of dimemas communication model for MPI collective operations *European PVM/MPI Users' Group Meeting* pp 39–46
- [12] Tikir M M, Laurenzano M, Carrington L and Snavely A 2009 PSINS: An open source event tracer and execution simulator for MPI applications *Euro-Par* (Delft, the Netherlands)
- [13] Alam S R and Vetter J S 2006 A framework to develop symbolic performance models of parallel applications *IPDPS* (IEEE)